

10 Game Theory

10.1 Games

Game theory studies decision problems in which the outcome of an agent's choice depends on other agents' choices. Such problems are called **games**, and the agents **players**. The Prisoner's Dilemma (example 1.3) is a game in this sense, because the outcome of your choice (confessing or remaining silent) depends on what your partner decides to do.

Whenever an agent faces a choice in a game, the MEU Principle tells us that they ought to choose whichever option maximizes expected utility. We don't need a new decision theory for games. Nonetheless, there are reasons for studying the special case where the states in a decision problem are other people's (real or potential) actions.

One reason is that we may be able to shed light on important social and political issues. The way we live and behave, as a society, is in many ways not ideal. We are depleting the Earth's resources. We are destabilising the climate. We are woefully underprepared for pandemics and other catastrophes. We buy goods from online retailers where most of the products are a scam. Corruption is rampant. The political system is broken. Dating is broken. And so on, and on. Why? Why don't we fix these problems? Is it because the current system benefits powerful actors who have us under their control? Game theory suggests an alternative possibility.

Remember the Prisoner's Dilemma. If you and your partner are rational and don't care about each other, you both confess and spend a long time in prison. Collectively, you could have achieved a much better outcome by remaining silent. Things are unnecessarily bad – you spend a long time in prison – not because a powerful third party stands to gain from your misery. The bad outcome is simply a result of your misaligned incentives.

This kind of situation is sadly common. Professional athletes, for example, have a strong incentive to use steroids, as long as the chance of being caught is low. Whether

or not their competitors do the same, using steroids provides an advantage. The outcome is that everyone uses steroids, even though everyone would prefer that no one uses steroids. Structurally, the athletes' decision problem is the same as the Prisoner's Dilemma. Any decision problem with this structure is nowadays called a Prisoner's Dilemma, even if no prisoners are evolved.

Another famous example is the "tragedy of the commons". Fishermen have an incentive to catch as many fish as they can, even though everyone would be better off if everyone restrained themselves to sustainable quotas.

Thomas Hobbes (in effect) argued that the pervasiveness of Prisoner's Dilemmas justifies the subordination of people under a state. It is in everyone's interest to impose a system of control and punishment that ensures the best outcome in what would otherwise be a Prisoner's Dilemma.

Exercise 10.1 †

Explain why a system of control and punishment can change a decision problem from a Prisoner's Dilemma into a problem with a different structure.

Another reason to study games is that a new set of conceptual tools and techniques become available if the states in a decision problem are other people's actions. In particular, we can often figure out which state obtains based on the other players' desires. In the (original) Prisoner's Dilemma, we know that if your partner is rational and only cares about their own prison term then they will confess.

Here is how game theorists would typically draw the matrix for the Prisoner's Dilemma, assuming you and your partner don't care about each other:

	Confess	Silent
Confess	-5, -5	0, -8
Silent	-8, 0	-1, -1

As before, the rows are the acts available to you. The columns are the acts available to your partner. We generally don't assign credences to the columns. The numbers in the cells represent the utility of the relevant outcome for you and your partner. We don't describe the outcome itself any more, for lack of space. The first number in each cell is the utility for the row player (whom we'll call 'Row' and assume to be female); the second is the utility for the column player ('Column', male).

In game theory jargon, a **solution** to a game is a prediction of what each player is going to do, assuming that they are rational. The solution to the Prisoner's Dilemma is that each player confesses. Confessing dominates remaining silent. You should confess no matter what you think your partner will do.

Consider the following matrix, for a different kind of game.

	C_1	C_2
R_1	2, 2	1, 3
R_2	1, 1	2, 2

Row no longer has a dominant option. What she should do depends on what she thinks Column will do. If Column chooses C_1 , then Row should play R_1 ; if Column chooses C_2 , then Row should play R_2 . Can we nonetheless say what Row will do, without specifying her beliefs?

Look at the game from Column's perspective. No matter what Row does, Column is better off choosing C_2 . C_2 dominates C_1 . So if Row knows the utility that Column assigns to the outcomes, then she can figure out that Column will choose C_2 . And so Row should choose R_2 . The solution is R_2, C_2 : Row chooses R_2 and Column C_2 .

Here is another, more complex example.

	C_1	C_2	C_3
R_1	0, 1	2, 2	3, 1
R_2	2, 2	1, 3	2, 2
R_3	1, 1	0, 2	0, 3

From Row's perspective, R_1 is the best choice if Column plays C_2 or C_3 , and R_2 is the best choice if Column goes for C_1 . For Column, C_2 is the best choice in case of R_1 or R_2 , and C_3 is best in case of R_3 . But Column can hardly expect Row to choose R_3 , since R_3 is dominated by R_2 . Column can figure out that Row will play either R_1 or R_2 , which means that Column will play C_2 . And since Row can figure out that Column will play C_2 , Row will play R_1 . The solution is R_1, C_2 .

To reach this conclusion, we need to assume more than that both players know each other's utilities. To figure out that Column will play C_2 , Row needs to know

that Column knows her (Row's) utilities, and she needs to know that Column knows that she (Row) won't choose a dominated option.

A common idealisation in game theory is that the players have **complete information** about the game, meaning that

- (1) all players know the structure of the game (as displayed in the matrix);
- (2) all players know that all players are rational;
- (3) all players know that (1)–(3) are satisfied.

By applying to itself, the clause (3) ensures that (1) and (2) hold with arbitrarily many iterations of 'all players know that' stacked in front. If something is in this way known by everyone, and known by everyone to be known by everyone, and so on, then it is said to be **common knowledge**. (1)–(3) say that the structure of the game and the rationality of all participants are common knowledge.

Exercise 10.2 ††

Under the assumptions (1)–(3), what will Row and Column do in the following games?

a.

	C_1	C_2
R_1	1, 0	1, 2
R_2	0, 3	0, 1

b.

	C_1	C_2	C_3
R_1	1, 0	1, 2	0, 1
R_2	0, 3	0, 1	2, 0

c.

	C_1	C_2	C_3
R_1	0, 1	2, 0	2, 4
R_2	4, 3	1, 4	2, 5
R_3	2, 4	3, 6	3, 1

10.2 Nash equilibria

Have a look at this game.

	C_1	C_2	C_3
R_1	4, 2	2, 3	3, 1
R_2	3, 1	3, 2	4, 1
R_3	4, 2	1, 1	0, 3

No option for either player is dominated by any other. Can we nonetheless figure out what Row and Column will choose?

Let's start with some trial and error. Take R_1, C_1 . Could this be the outcome that is reached whenever the game is played by two players under the idealizing assumptions (1)–(3)? No. Otherwise Column would know that Row is going to play R_1 . And then Column is better off playing C_2 . The opposite happens with R_1, C_2 : if Row knew that Column plays C_2 , she would be better off playing R_2 . This kind of reasoning disqualifies all combinations except R_2, C_2 – the middle cell. If Row knows that Column is going to play C_2 , she can do no better than play R_2 . Likewise for Column: if Column knows that Row is going to play R_2 , he can do no better than play C_2 .

A combination of options that is “stable” in this way is called a **Nash equilibrium** (after the economist John Nash). In general, a Nash equilibrium is a combination of acts, one for each player, such that no player could get greater utility by deviating from their part of the equilibrium, given that the other players stick to their part.

There is a simple algorithm for finding Nash equilibria in finite two-player games. Start from the perspective of the row player. For each act of the column player, underline the best outcome(s) Row can achieve if Column chooses this act. In the example above, you would underline the 4s in the first column, the 3 in the middle cell, and the 4 in the third column. Then do the same for the column player: for each act of Row, underline the best possible outcome(s) for Column. The result looks like this.

	C_1	C_2	C_3
R_1	<u>4</u> , 2	2, <u>3</u>	3, 1
R_2	3, 1	<u>3</u> , <u>2</u>	<u>4</u> , 1
R_3	<u>4</u> , 2	1, 1	0, <u>3</u>

Any cell in which both numbers are underlined identifies a Nash equilibrium.

If a game has a unique Nash equilibrium, and assumptions (1)–(3) hold, then the equilibrium is plausibly the game's solution: each player can be expected to play their part of the equilibrium.

This is not as obvious as it may perhaps appear. Consider the next game.

	C_1	C_2	C_3
R_1	<u>2</u> , -2	-1, <u>1</u>	<u>1</u> , -1
R_2	0, 0	<u>0</u> , 0	-2, <u>2</u>
R_3	0, <u>0</u>	<u>0</u> , <u>0</u>	<u>1</u> , -1

There is a unique Nash equilibrium: R_3, C_2 . If this is the guaranteed outcome under assumptions (1)–(3), then Row can be sure that Column will play C_2 . But if Column plays C_2 , then R_2 and R_3 are equally good for Row. So how can we be sure Row won't play R_2 ?

You might argue that if Row played R_2 and Column could predict her choice, then Column would play C_3 , leading to a worse result for Row. But we're not assuming that Column can predict Row's choice. All we're assuming is (1)–(3).

Still, there is an argument in favour of R_3, C_2 as the unique solution. Suppose for reductio that Row could play either R_3 or R_2 , and conditions (1)–(3) are satisfied. Then Column could not be sure which of these Row will choose; without further information, he would have to give roughly equal credence to R_2 and R_3 . But then his best choice is C_3 . Anticipating this, Row ought to choose R_3 . This contradicts our assumption that Row could just as well play R_2 .

Exercise 10.3 †

Identify the Nash equilibria in the following games.

a.

	C_1	C_2
R_1	3, 4	4, 3
R_2	1, 3	5, 2
R_3	2, 0	1, 5

b.

	C_1	C_2	C_3
R_1	1, 0	1, 2	0, 1
R_2	0, 3	0, 1	2, 0

c.

	C_1	C_2	C_3
R_1	0, 1	2, 0	2, 4
R_2	4, 3	1, 4	2, 5
R_3	2, 4	3, 6	3, 1

Exercise 10.4 ††

Whenever the method from section 10.1, which is called **elimination of dominated strategies**, identifies a combination of acts as a game's solution, then this combination of acts is a Nash equilibrium. Can you explain why?

10.3 Zero-sum games

In some games, the players' preferences are exactly opposed: if Row prefers one outcome to another by a certain amount, then Column prefers the second outcome to the first by the same amount. The utilities in every cell sum to the same number. Since utility scales don't have a fixed zero, we can re-scale the utilities so that the sum is zero. For this reason, games in which the players' preferences are opposed are called **zero-sum games**. Here is an example.

	C_1	C_2	C_3
R_1	1, <u>-1</u>	<u>3</u> , -3	<u>1</u> , <u>-1</u>
R_2	<u>2</u> , -2	-2, <u>2</u>	-1, 1

There is a unique Nash equilibrium: R_1, C_3 . Under assumptions (1)–(3), the MEU Principle entails that Row should play R_1 and Column C_3 – although this isn't obvious, since we haven't specified any probabilities. Curiously, R_1, C_3 is also supported by the *maximin* rule that we've met in section 1.4. Maximin tells each player to choose an option with the best worst-case result. In our example, the worst-case result of choosing R_1 (for Row) has utility 1; the worst-case result of R_2 is -2. Maximin therefore says that Row should choose R_1 . For Column, it similarly recommends C_3 .

This is not a coincidence. Every Nash equilibrium in every zero-sum game is supported by the maximin rule. For suppose it isn't. Suppose, more concretely, that R_i, C_j is a Nash equilibrium in a (two-player) zero-sum game, but R_i isn't supported by the maximin rule. Then the outcome of R_i, C_j isn't the worst possible outcome of R_i for Row: there is some alternative C_k to C_j for which R_i, C_k is worse for Row than R_i, C_j . Since the game is zero-sum, R_i, C_k is *better for Column* than R_i, C_j . And so R_i, C_j isn't a Nash equilibrium.

Many games have more than one Nash equilibrium. As we will see, it can then be hard to predict what the players will do without looking at their beliefs. They may not even reach one of the Nash equilibria. In zero-sum games, however, this problem is unlikely to arise. Consider the following example.

	C_1	C_2	C_3
R_1	2, -2	<u>1</u> , <u>-1</u>	<u>1</u> , <u>-1</u>
R_2	<u>3</u> , -3	<u>1</u> , <u>-1</u>	<u>1</u> , <u>-1</u>
R_3	0, 0	-1, 1	-2, <u>2</u>

The game has four Nash equilibria. What will the players do? Should Row play R_1 or R_2 ? Should Column play C_2 or C_3 ? Well, it doesn't matter. The players can arbitrarily choose among these options. Whatever they choose, they are guaranteed to end up at an equilibrium, and all the equilibria have the same utility.

Exercise 10.5 †††

Prove that this holds for all two-player zero-sum games: if R_i, C_j and R_n, C_m are Nash equilibria, then so are R_i, C_m and R_n, C_j ; moreover, all Nash equilibria have the same utility.

Some games have no Nash equilibrium at all. Here is a matrix for Rock–Paper–Scissors.

	Rock	Paper	Scissors
Rock	0, 0	-1, <u>1</u>	<u>1</u> , -1
Paper	<u>1</u> , -1	0, 0	-1, <u>1</u>
Scissors	-1, <u>1</u>	<u>1</u> , -1	0, 0

There is no equilibrium. What should you do in this kind of game?

A standard answer in game theory is that you should randomize. You should, say, toss a fair die and choose Rock on 1 or 2, Paper on 3 or 4, and Scissors on 5 or 6. Such a randomized choice is called a **mixed strategy**. We will write '[1/3 Rock, 1/3 Paper, 1/3 Scissors]' for the mixed strategy of playing Rock, Paper, or Scissors each with (objective) probability 1/3.

Suppose two players both play [1/3 Rock, 1/3 Paper, 1/3 Scissors]. Then neither could do better by playing anything else (including other mixed strategies). The combination of the two mixed strategies is a Nash Equilibrium. It is the only Nash Equilibrium in Rock–Paper–Scissors.

It can be shown that every finite game has at least one Nash Equilibrium if mixed strategies are included. (This was shown by John Nash.) The proof obviously assumes that randomization introduces no additional costs or benefits. If you hate randomization and prefer losing in Rock–Paper–Scissors to randomizing, then the game has no Nash Equilibrium, not even among mixed strategies.

Exercise 10.6 ††

Suppose your opponent plays $[1/3 \text{ Rock}, 1/3 \text{ Paper}, 1/3 \text{ Scissors}]$. What is the expected utility of playing Rock? How about Paper and Scissors? What is the expected utility of playing $[1/3 \text{ Rock}, 1/3 \text{ Paper}, 1/3 \text{ Scissors}]$?

10.4 Harder games

Most games in real life are not zero-sum games. The following example illustrates the class of **coordination problems** in which the players would like to coordinate their actions.

Example 10.1

You and your friend Bob want to meet up, but neither of you knows to which party the other will go. Party A is better than party B, but you will both go home if you don't find each other.

	Party A	Party B
Party A	3, 3	0, 0
Party B	0, 0	2, 2

There are two Nash equilibria (without randomization): both going to party A, and both going to party B. We can't assume that whenever rational agents play the game, then they will end up in one of these equilibria. If you suspect that Bob will go to party B, and Bob suspects you will go to party A, then you'll go to B and Bob to A.

But could this actually happen, under assumptions (1)–(3)? As you may check, going to party B maximizes expected utility if and only if your credence that Bob goes to B is at least 0.6. But could you be at least 60% confident that Bob will go to B, given what you know about Bob’s utilities?

Well, Bob will go to B provided that *he* is at least 60% confident that *you* will go to B. So to be at least 60% confident that Bob will go to B, you only need to be at least 60% confident that Bob is at least 60% confident that you will go to B. Of course, Bob can figure out that you will go to B only if you are at least 60% confident that he will go to B. So to be at least 60% confident that Bob will go to B, you need to be at least 60% confident that Bob is at least 60% confident that you are at least 60% confident that Bob will go to B. And so on. There is nothing incoherent about this state of mind, in which you are at least 60% confident that Bob will go to B. Nonetheless, we may wonder how you could have arrived at it. How could you have rationally arrived at a 60% confidence that Bob is at least 60% confident that you are at least 60% confident that ...and so on and on forever?

The assumptions (1)–(3) here give rise to an epistemological puzzle. If you have no further relevant evidence, how confident should you be that Bob goes to B? You might think your degree of belief should be $1/2$, by the Principle of Indifference. But then you should assume that Bob’s degree of belief in *you* going to B is also $1/2$. And that would imply that Bob goes to A. So it can’t be right that you should give equal credence to the two possibilities.

In real coordination problems, the players often do have further information. When you’re driving on a road, you are playing a coordination game with drivers going in the opposite direction. You prefer to drive on the left if and only if the others drive on the left; the others prefer to drive on the left if and only if you drive on the left. The existence of a law to drive on the left gives you reason to think that the others will drive on the left. But even without a law, the mere observation that people generally drive on the left would give you reason to think that that’s what they will continue to do.

A different kind of coordination is called for in the following game.

Example 10.2 (Chicken)

For fun, you and your friend Bob drive towards each other at high speed. If one of you swerves and the other doesn't, the one who swerves loses. If neither swerves, you both die.

	Swerve	Straight
Swerve	0, 0	-1, 1
Straight	1, -1	-10, -10

Games like chicken are sometimes called **anti-coordination games**, because each player would prefer the other one to yield without yielding themselves. There are two Nash Equilibria in Chicken that don't involve randomization: 'Swerve, Straight' and 'Straight, Swerve'. As above, either choice is rationally defensible, given suitable beliefs about the opponent, and as before there is an epistemological puzzle about how any of these beliefs could come about.

An interesting feature of many anti-coordination games is that they seem to favour irrational agents who do not maximize expected utility. Suppose Bob is insane and will go straight no matter what, despite the large cost of dying if you both go straight. And suppose you know about Bob's insanity. Then you, as an expected utility maximizer, will have to swerve. Bob will win.

There are stories that during the cold war, the CIA leaked false information to the Russians that the US President was an alcoholic, while the KGB falsified medical reports suggesting that Brezhnev was senile. Both sides tried to gain a strategic advantage over the other by indicating that they would irrationally retaliate against a nuclear strike even if they had nothing to gain any more.

Exercise 10.7 †

What should you do in Chicken if you give equal credence to the hypotheses that Bob will swerve and that he will go straight?

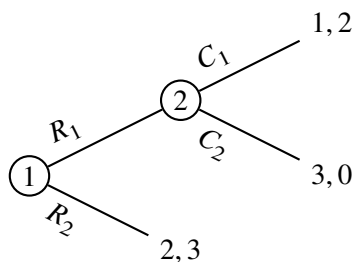
Exercise 10.8 †††

A third Nash equilibrium in Chicken involves randomization. Can you find it? What is the expected utility for both players if they play the mixed strategy?

10.5 Games with several moves

So far, we have looked at games in which each player makes just one move, and no player knows about the others' moves ahead of their choice. Game theory also studies situations in which these assumptions are relaxed. Let's have a quick look at games with several moves, assuming players always know what was played before.

As in section 8.2, we can picture the relevant decision situations in a tree-like diagram (an "extensive form representation"). Below is a diagram for a game in which Row first has a choice between R_1 and R_2 . If she chooses R_2 , the game ends with an outcome that has utility 2 for Row and 3 for Column. If Row chooses R_1 , then Column gets a choice between C_1 and C_2 . If he chooses C_2 , Row gets utility 3 and Column 0; if Column chooses C_1 , Row gets 1 and Column 2.



We can use backward induction to predict how the game is going to be played, assuming (1)–(3).

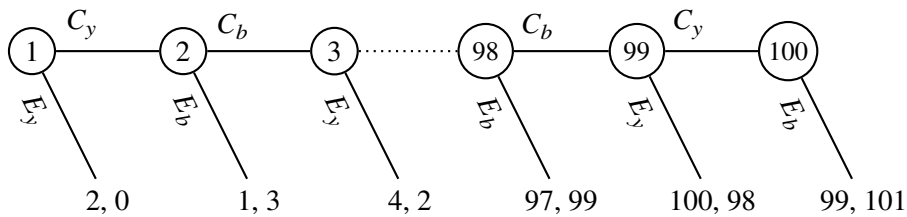
Consider node 2, where Column has a choice between outcome '3, 0' and outcome '1, 2'. The choice involves no relevant uncertainty, and Column prefers '1, 2' over '3, 0'. He can be expected to play C_1 . Anticipating this, Row can figure out that playing R_1 at node 1 will lead to '1, 2'. R_1 instead leads to '2, 3'. This is better for Row. So Row will play R_1 .

In the following example, backward induction leads to a more surprising result.

Example 10.3 (Centipede)

You and Bob are playing a game. The game starts with a pot containing £2. In round 1, you can decide whether to continue or end the game. If you end the game, you get the £2 and Bob gets £0. If you continue, the money in the pot increases by £2 and Bob decides whether to continue or end. If he ends the game here (in round 2), the pot is divided so that he gets £3 and you get £1. If he continues, the money in the pot increases by another £2 and it's your turn again. If you end the game (in round 3), you get £4 and Bob gets £2. And so on. In each round, the money in the pot increases by £2 and whoever ends the game gets £2 more than the other player. In round 100, Bob no longer has an option to continue.

Suppose you and Bob don't care about each other; each of you only wants to get as much money as possible. Here is a partial diagram of the resulting game.



Let's use backward induction to solve the game. At node 100, Bob doesn't have a choice. If you continue at node 99 (C_y), you will get £99 and Bob £101. If you end the game (E_y) at node 99, you will get £100. It is obviously better to end the game. Anticipating this, what should Bob do in round 98? If he ends the game (E_b), he'll get £99; if he continues (C_b), he'll get £98. So he should end the game. Anticipating this, you should end the game in round 97, to ensure that you'll get £98 rather than £97. And so on, all the way back to round 1. At each point, backward induction tells us that the game should be ended. In particular, you can anticipate in round 1 that Bob will end the game in round 2. So you should end the game in round 1. You will get £2 and Bob £0.

When actual people play the Centipede game, almost no-one ends the game right away. Is this a sign of either altruism or irrationality? Not necessarily.

Let's look at your choice in round 1 from an MEU perspective. It is clear what happens if you end the game: you'll get £2. But what would happen if you chose to continue? The argument from backward induction assumes that Bob would end the game. If you could be certain that Bob would do that, then you should indeed end the game in round 1. But why should Bob end the game? Because, so the argument, he can be certain that you would end the game in round 3. But the argument for ending in round 3 is exactly parallel to the argument for ending in round 1. And if Bob faces a choice in round 2, then he has just seen that you *did not* end the game in round 1. Based on this information, he can't be sure you would end it in round 3. On the contrary, he should be somewhat confident that you will continue in round 3. And then continuing maximizes expected utility in round 2. Anticipating this, continuing also maximizes expected utility in round 1, as it is likely to get you at least to round 3.

This suggests that the backward induction argument went wrong somewhere. But where? Surely you really ought to end the game in round 99. And surely this means that Bob should end the game in round 98. And so on! This puzzle is sometimes called the **paradox of backward induction**.

Exercise 10.9 ††

Consider a variant of the Centipede game with no fixed end point. Instead, each time a player chooses to continue, the game ends with a probability of 1%. Does this change anything? How should you play?

Exercise 10.10 ††

Suppose you repeatedly face the Prisoner's Dilemma with the same partner, for an unknown number of rounds. You only care about your own prison terms. You expect that your partner will remain silent in the first round and from then on imitate whatever you did in the previous round. What should you do? Does your answer show that you should choose a dominated act?

10.6 Evolutionary game theory

One of the most successful applications of game theory lies (somewhat surprisingly) in the study of biological and cultural evolution. Consider the following game.

Example 10.4 (The Stag Hunt)

Two players independently decide whether to hunt stag or rabbit. Hunting stag requires cooperation, so if only one of the players decides to hunt stag, she will get nothing. The utilities are as follows.

	Stag	Rabbit
Stag	5, 5	0, 1
Rabbit	1, 0	1, 1

In the evolutionary interpretation, the utilities represent the *relative fitness* that results from a combination of choices, measured in terms of average number of surviving offspring. Let's assume that each strategy is played by a certain fraction of individuals in a population. Individuals who achieve an outcome with greater utility will, by definition, have more offspring on average, so their proportion in the population will increase.

Suppose initially $1/4$ of the individuals in the population goes for stags and $3/4$ for rabbits. Assuming that encounters between individuals are completely random, this means that any given individual has a $1/4$ chance of playing with someone hunting stag, and a $3/4$ chance of playing with someone hunting rabbit. The average utility of hunting stag is $1/4 \cdot 5 + 3/4 \cdot 0 = 1.25$; for hunting rabbit the utility is of course 1. Individuals going for stag have greater average fitness. Their fraction in the population increases. As a consequence, it becomes even more advantageous to go for stag. Eventually, everyone will hunt stag.

By contrast, suppose initially only $1/10$ of the population goes for stags. Then hunting stag has an average utility of 0.5, which is less than the utility of hunting rabbit. The rabbit hunters will have more offspring, which makes it even worse to hunt stags. Eventually, everyone will hunt rabbits.

The two outcomes ‘Stag, Stag’ and ‘Rabbit, Rabbit’ are the two Nash Equilibria in the Stag Hunt. Evolutionary game theory predicts that the proportion of stag and rabbit hunters in a population will approach one of these equilibria.

Not every Nash Equilibrium is a possible end point of evolution though. If a population repeatedly plays the game of Chicken, and the players can’t recognize in advance who will swerve and who will go straight, then the asymmetric equilibria ‘Swerve, Straight’ and ‘Straight, Swerve’ do not mark possible end points of evolutionary dynamics. But note that in a community in which almost everyone swerves, you’re better off going straight; similarly, in a community in which almost everyone goes straight, the best choice is to swerve. Evolution will therefore lead to the third, mixed strategy equilibrium. It will lead to a state in which a certain fraction of the population swerves and the others go straight.

The assumption that individuals in a population are randomly paired with one another is obviously an idealisation. In reality, individuals are more likely to interact with members of their own family, which increases the chances that they will be paired with individuals of the same type; they might also actively seek out others who share the relevant traits. Either way, the resulting **correlated play** dramatically changes the picture.

Imagine a population in which individuals repeatedly play a Prisoner’s Dilemma wherein they can either cooperate (remain silent, in the original scenario) or defect (confess). Since defectors do better than cooperators in any encounter, it may seem that cooperation can never evolve. On the other hand, cooperators do much better when paired with other cooperators than defectors when paired with defectors. If the extent of correlation is sufficiently high, cooperators can take over (although perhaps not completely).

In many species, one can find altruistic individuals who sacrifice their own fitness for the sake of others. Evolutionary game theory explains how this kind of altruism could have evolved.

Exercise 10.11 †

Why can’t we expect cooperative behaviour to take over completely in the scenario where cooperation spreads through correlated play?

Exercise 10.12 †

What are the Nash equilibria in the following game (ignoring randomization)? Could all the equilibria come about through an evolutionary process?

	A	B
A	5, 5	1, 1
B	1, 1	1, 1

Essay Question 10.1

Explain the paradox of backward induction. Why is it a paradox? How do you think it could be resolved?

Sources and Further Reading

There are many decent introductions to Game Theory. The “[Game Theory](#)” entry in the Stanford Encyclopedia by Don Ross (2019) provides a fairly comprehensive overview. A suitable next step might be Steven Tadelis, *Game Theory: An Introduction* (2013).

The paradox of backward induction is discussed, for example, in Philip Pettit and Robert Sugden, “The Backward Induction Paradox” (1989).

For a little more on evolutionary game theory, see Brian Skyrms, “Game Theory, Rationality and Evolution of the Social Contract” (2000). For even more, see Brian Skyrms, “The Stag Hunt and the Evolution of Social Structure” (2004).